

1 Representation Ensemble Learning Applied to
2 Facial Expression Recognition

3 Bruna Rossetto Delazeri^{1*†}, Andre Gustavo Hochuli^{1†},
4 Jean Paul Barddal^{1†}, Alessandro Lameiras Koerich^{2†},
5 Alceu de Souza Britto Jr.^{1,3†}

6 ¹Graduate Program in Informatics (PPGIA), Pontifícia Universidade
7 Católica do Paraná (PUCPR), Imaculada Conceição, Curitiba,
8 80215-901, Paraná, Brazil.

9 ²Department of Software and IT Engineering, École de Technologie
10 Supérieure (ÉTS) - University of Québec, 1100 Notre-Dame St W,
11 Montreal, H3C 1K3, Quebec, Canada.

12 ³Postgraduate Program in Applied Computing, State University of
13 Ponta Grossa (UEPG), Av. General Carlos Cavalcanti, Ponta Grossa,
14 84010-330, Paraná, Brazil.

15 *Corresponding author(s). E-mail(s): bruna.delazeri@pucpr.edu.br;

16 Contributing authors: aghochuli@ppgia.pucpr.br;

17 jean.barddal@ppgia.pucpr.br; alessandro.koerich@etsmtl.ca;

18 alceu@ppgia.pucpr.br;

19 †These authors contributed equally to this work.

20 **Abstract**

21 This work introduces the representation ensemble learning algorithm, a novel
22 approach for generating diverse unsupervised representations rooted in the prin-
23 ciples of self-taught learning. The ensemble comprises convolutional autoencoders
24 (CAEs) learned in an unsupervised manner, fostering diversity via a loss function
25 designed to penalize similar CAEs' latent representations. We employ support
26 vector machines, bagging, and random forest as primary classification methods
27 for the final classification step. Additionally, we incorporate KnoraU, a well-
28 established technique used to dynamically select competent classifiers based on a
29 test sample. We evaluate various fusion strategies, including sum, product, and
30 stacking, to comprehensively assess the ensemble's performance. A robust experi-
31 mental protocol considering the facial expression recognition problem shows that

32 the proposed approach based on self-taught learning surpasses the accuracy of
33 fine-tuned convolutional neural network (CNN) models. In terms of accuracy,
34 the proposed method is up to 9.9 and 6.3 percentage points better than the
35 CNN-based models fine-tuned for JAFFE and CK+ datasets, respectively.

36 **Keywords:** Self-taught Learning, Facial Expression Recognition, Convolutional
37 Neural Network, Autoencoder

38 1 Introduction

39 The study of facial expression recognition (FER) has been ongoing for over two
40 decades, as it conveys emotional states during human communication that play an
41 essential role in society. Moreover, recognizing emotions has several applications, some
42 of them with substantial social impact, such as identifying the autism spectrum [1] and
43 chronic depression disorders [2], or monitoring driver fatigue in safety car systems [3].

44 The main challenges on the FER task are related to the different environmen-
45 tal conditions like lighting variations, non-frontal image acquisition, low-intensity face
46 expression, and differences in facial expressions among gender, culture, and age groups
47 [4]. Additionally, the scarcity of training data concerning its quantity and quality
48 (samples of facial emotions) for some specific FER applications makes a full supervised-
49 based approach bothersome [5]. In other words, a significant limitation of purely
50 supervised systems is their dependence on large volumes of labeled data. Labeling is,
51 in some FER applications, a costly and time-consuming process. The deep learning-
52 based methods disseminated with convolutional neural network (CNN) models have
53 made such a limitation even more evident [6].

54 An up-and-coming alternative to mitigate the problem related to the lack of labeled
55 data is transfer learning (TL). A particular case of TL, named self-taught learning
56 (STL) [7], consists of learning a representation (feature extractor) from unlabeled data,
57 which is later applied to a target domain that is potentially from a different distri-
58 bution. The main idea behind STL relies on natural human learning where unlabeled
59 data is considered an essential foundation for high-level learning, being responsible for
60 providing a more significant discriminating power [7]. Even with enough labeled data,
61 STL can enhance the learning process, as observed in the experiments conducted in
62 [8].

63 The combination of extraction and classification techniques is essential for FER
64 systems. The main challenges in efficient facial expression recognition are the selection
65 of efficient feature extraction and classification techniques. Handcrafted features, such
66 as local binary patterns (LBP) [9, 10], scale-invariant feature transforms (SIFT) [11],
67 gray level co-occurrence matrix (GLCM) [12], sped up robust features (SURF) [13]
68 and histograms of oriented gradient (HOG) [14] achieved a breakthrough in various
69 fields. Although using different descriptors to extract handcrafted features is possible,
70 it is currently common to choose to learn the representation, removing the task of
71 defining the features from the developer.

72 The potential of deep learning approaches to generate highly effective data repre-
73 sentations has been widely recognized. In FER research, solutions based on CNNs have
74 gained prominence [15]. Similarly, the convolutional autoencoder (CAE) is a promising
75 strategy to make unsupervised-based solutions [16, 17] possible. In our previous study
76 [18], we noticed that unsupervised learned representations yielded an effective facial
77 emotion recognition method. Built upon such prior work, we introduce in this paper
78 a new representation learning ensemble (REL) algorithm to produce a pool of repre-
79 sentations employing STL principles. The main idea behind REL is to explore various
80 strategies to learn representations (feature extractors) unsupervised and consider a
81 custom loss function designed to promote diversity. Moreover, in the context of the
82 FER problem, we evaluate how competitive the REL results are compared to CNN-
83 based architectures trained in a supervised approach. We also compare our method
84 with the current state-of-the-art in a common experimental protocol.

85 Two research questions guide this paper. The first, referred to as RQ1, concerns
86 generating a diverse pool of unsupervised representations: “Does using a pool of
87 unsupervised representations generated by the REL algorithm contribute to the per-
88 formance of FER solutions?” However, we believe that more important than providing
89 a new STL-based method to generate a pool of complementary representations is
90 showing how far its performance is from that of supervised-based solutions. Thus,
91 our second research question (RQ2) is “How does the new REL algorithm com-
92 pare to supervised CNN-based architectures regarding facial expression recognition
93 performance?”

94 By answering these research questions, the present manuscript provides a two-
95 fold contribution. The first contribution is a new algorithm to generate unsupervised
96 representations using STL focusing on diversity. More specifically, CAEs are trained
97 considering different diversity induction strategies, i.e., distinct model initialization,
98 architectures, and training data. The rationale is to investigate if complementary un-
99 supervised learned representations can contribute to the performance of a FER solution.
100 The second contribution is a robust performance comparison between the proposed
101 unsupervised FER solution and supervised strategies under the same conditions, i.e.,
102 the same experimental protocol, and considering an authentic and demanding task
103 like FER.

104 A robust experimental protocol considering three auxiliary datasets (Kyoto, LFW,
105 and LabelMe) necessary for the unsupervised generation of representations proposed
106 here and two well-known FER datasets (Jaffe and CK+) has shown that the proposed
107 method may surpass the accuracy of CNN-based solutions, including two fine-tuning
108 strategies applied to ten different CNN architectures and the training of two CNNs
109 from scratch. Furthermore, the proposed FER framework demonstrated superiority
110 over several studies recently published in high-impact journals, covering both deep
111 learning approaches and handcrafted frameworks and following the leave-one-subject-
112 out (LOSO) cross-validation protocol.

113 This paper is organized as follows. Section 2 introduces related works. Section
114 3 describes the proposed algorithm for learning a pool of diverse representations,
115 while Section 4 presents the experiments undertaken to answer our research questions.
116 Finally, Section 5 draws our conclusions and directions for future work.

117 2 Related Works

118 Various techniques have been proposed as FER solutions in the last few decades.
119 Ensemble learning has become widely utilized in emotion recognition owing to its
120 superior accuracy and generalization. Exploring diversity to create ensembles of CNNs,
121 Renda et al. [19] investigated different strategies for inducing diversity in an ensemble
122 of CNNs applied to FER. The results on the FER2013 dataset showed that seed
123 variation yielded the best recognition results, while variations on the pre-training
124 process of their CNNs achieved the best run-time performance.

125 A novel method for FER using an ensemble of CNNs with probability-based fusion
126 was proposed by Wen et al. [20]. The ensemble is constructed by randomly varying
127 parameters and architecture around the optimal values for CNN. Each CNN is trained
128 to output a probability for each class, which is then fused using probability-based
129 fusion. Experiments were conducted using 100 CNNs, and the best average results
130 were achieved when the ensemble size was 35 for all datasets. The method was tested
131 on benchmark datasets, achieved 50.70% and 76.05% for JAFFE and CK+ datasets,
132 respectively, and outperformed other compared methods in terms of accuracy.

133 An ensemble learning method using electroencephalogram (EEG) signals is pro-
134 posed by Li et al. [21]. The method uses a sliding time window to extract features and
135 L1 regularization to select effective features. It then applies a model selection method
136 to choose the best basic analysis sub-models and an ensemble operator to convert
137 classification results. The optimal parameters are determined using multiple objective
138 particle swarm optimization. The method is evaluated on two public datasets (DEAP
139 and SEED) using the LOSO experimental protocol. The average accuracy rates for
140 arousal and valence are 65.70% and 64.22%, respectively, on the DEAP dataset, and
141 the average accuracy on the SEED dataset is 84.44%.

142 TL has gained popularity and emerged as a promising area in machine learning due
143 to its broad potential applications and has been extensively applied to computer vision
144 tasks. Dhankhar et al. [22] introduced ResNet50 and VGG16 architectures for facial
145 emotion recognition and proposed an ensemble that combines both CNN models. The
146 ensemble model outperformed the baseline SVM as well as the individual ResNet50
147 and VGG16 networks, achieving the highest accuracy of 75.8%. The SVM had an
148 accuracy of 37.9%, while ResNet50 and VGG16 had accuracy rates of 73.8% and
149 71.4%, respectively. The authors also explored TL for FER using pre-trained AlexNet,
150 VGG, and ResNet architecture networks. They achieved an average accuracy of 90%
151 on the combined JAFFE and CK+ datasets.

152 Chowdary et al. [23] discuss using TL techniques in the FER scenario. Pre-trained
153 networks such as ResNet50, VGG19, Inception V3, and MobileNet are employed. The
154 fully connected layers of these pre-trained CNNs are removed and replaced with their
155 own fully connected layers suitable to the specific task requirements. These newly
156 added layers are then trained to update the weights. The experiment was performed
157 using the CK+ dataset, resulting in an average accuracy of 96%.

158 A novel pipeline strategy that gradually improves the accuracy of FER by first
159 training the dense layer(s) and then tuning each pre-trained CNN block successively
160 was proposed by Akhand et al. [24]. The proposed FER system is tested on eight
161 pre-trained CNN models using the KDEF and JAFFE facial image datasets. FER is

162 challenging even for frontal views and is further complicated by the diversity of pro-
163 file views in the KDEF dataset. The proposed method achieved remarkable accuracy
164 on both datasets with pre-trained models, achieving a FER accuracy of 96.51% and
165 99.52% on KDEF and JAFFE datasets, respectively, on a 10-fold cross-validation way.

166 Proposed by Raina et al. [7], the STL framework is a particular case of TL that uses
167 unlabeled data from distributions other than of the problem at hand to learn a high-
168 level representation in an unsupervised manner. This approach is motivated by the
169 lack of training instances in many applications, and authors argue that a robust rep-
170 resentation can be obtained from unlabeled samples, such as low-level discriminating
171 structures in natural images, such as corners, curves, and shapes [25]. As a conse-
172 quence, several authors have employed STL in diverse classification tasks, including
173 audio [26], text [27, 28], image [29–31] and sensor data [32].

174 Inspired by STL to mitigate the problem related to the lack of samples for training
175 in some FER applications and the use of ensembles to provide a way to generate a
176 pool of diverse representations, we present our proposed algorithm algorithm in the
177 next section.

178 3 Proposed Approach

179 The novelty in the proposed approach to solving the FER problem is an algorithm
180 to generate a pool of representations inspired by STL and diversity concepts. Our
181 approach is called representation ensemble learning (REL). REL has two essential
182 features to promote diversity: i) it explores different strategies to generate diverse CAE
183 architectures; ii) it uses a custom loss function to train different CAEs and promote
184 diversity amongst latent representations.

185 Fig. 1 presents a general overview of the complete STL strategy adopted in this
186 paper. In the **unsupervised representation learning step (Step 1)**, a high-level
187 representation is learned using unlabeled data. This is the part in which our proposal
188 is novel. Next, in the **feature building step (Step 2)**, feature vectors are extracted
189 from the labeled data of the target domain using the representation learned in the
190 previous step. Finally, in the **supervised learning step (Step 3)**, the feature vec-
191 tors extracted in Step 2 are used to train a FER classification model (either using
192 monolithic or ensemble-based classifiers).

193 3.1 Strategies Explored to Generate Diversity

194 We observe in the literature robust evidence that greater diversity is highly correlated
195 with the increase in supervised CNN-based ensemble accuracy [19, 33]. Consequently,
196 we follow the rationale of inducing diversity when creating the pool of unsupervised
197 representations as illustrated in Fig. 2. The proposed algorithm uses a CAE to auto-
198 matically learn meaningful representations from raw data, eliminating the need for
199 manual feature engineering. The different strategies to vary specific CAE parameters
200 are as follows:

- 201 • **Random Seed (S)**: here, different representations are generated by varying the
202 distribution of weights during the CAE initialization process. The architecture may

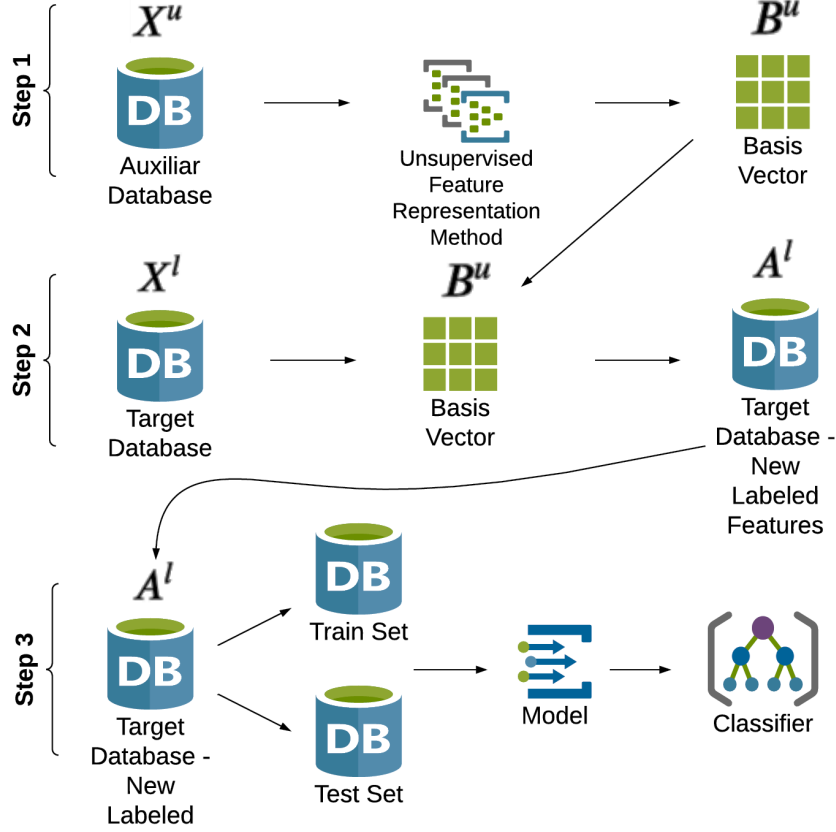


Fig. 1 STL steps: (1) Unsupervised representation learning; (2) Feature building; (3) Supervised learning.

203 be the same from one CAE to another, but the seed of the initialization process
 204 differs. The input is the number of representations (R) generated with different
 205 seeds. The algorithm randomly selects a seed at each generated CAE in the $[0, 1000]$
 206 interval.

- 207 • **CAE Architecture (A):** here, we explore representation generation using differ-
 208 ent CAE architectures. One must define the network depth (number of layers, D),
 209 the filters used, and the dimensionality (I) of the intermediate (latent) layer. The
 210 generator will create the first architecture with a depth equal to the number of
 211 defined layers (D), then the second will have $(D - 1)$ up to the last architecture
 212 with depth $D = 1$, which has one input, one intermediate, and one output layer
 213 (basic structure of a CAE). The representations generated here differ in depth on
 214 the CAE architectures but use the same number of neurons (I) in the latent layer.
- 215 • **Latent Representation (L):** in this strategy, diversity is induced by varying the
 216 number of neurons in the middle layer of the CAE, named the latent layer. The

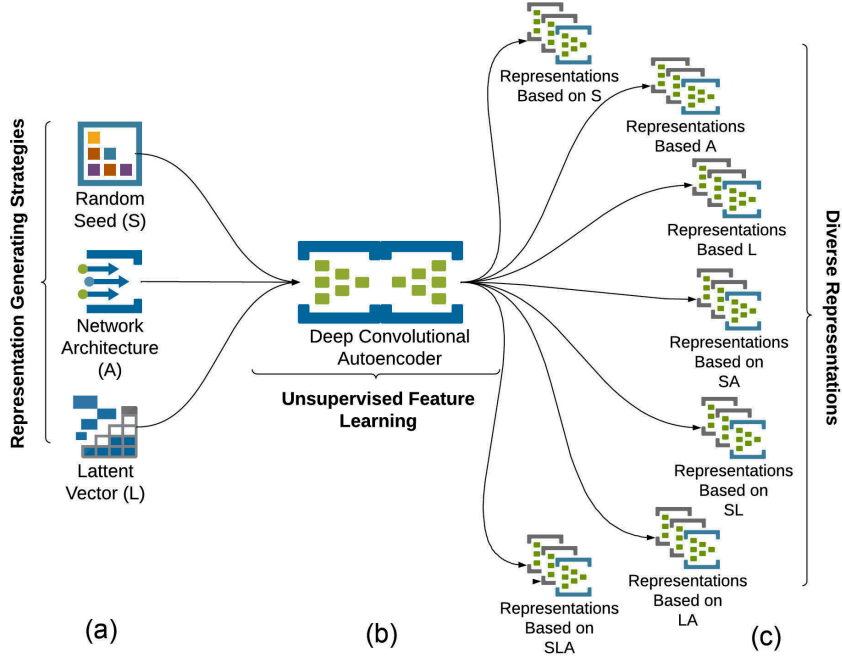


Fig. 2 Automatic unsupervised strategy generator for problem representation. (a) Different generating strategies of representations. (b) CAE for learning high-level representations. (c) Generated representations based on strategies: seeds (S), CAE architecture (A), latent representation (L), seeds + CAE architecture (SA), seeds + latent representation (SL), latent representation + CAE architecture (LA), and seeds + latent representation + CAE architecture (SLA).

217 input is the number of representations (R) to be generated with different dimen-
 218 sionality for the latent layer. The algorithm randomly selects a dimensionality from
 219 the $[150, 2500]$ range at each iteration.

220 The aforementioned strategies to generate diversity can also be combined as
 221 follows:

- 222 • **SA**: it combines the strategy based on seeds (S) and CAE architecture (A). The
 223 representations generated by this combination have variations in CAE initialization
 224 and architecture (number of layers).
- 225 • **SL**: it explores the combination of seeds (S) and latent representation (L) strategies,
 226 where the variation of both CAE initialization and latent layer dimensionality are
 227 adopted for each representation generated.
- 228 • **LA**: it explores the use of latent representation (L) and CAE architecture (A).
 229 Here, we use two approaches directly related to the CAE architecture, varying
 230 the dimensionality of the latent layer and the CAE depth in each representation
 231 generated.
- 232 • **SLA**: in the last combination, we consider all the strategies seeds (S), latent
 233 representation (L), and CAE architecture (A) for each representation generated.

234 **3.2 Custom Loss Function**

235 The loss function used in the REL algorithm contains two terms: a general loss term
 236 and a penalty term. The first term is the mean squared error (MSE) denoted in
 237 Equation 1:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{P} \sum_{i=1}^P (X_i - \hat{X}_i)^2 \quad (1)$$

238 where X_i is the original input data, \hat{X}_i is the output data reconstructed by the CAE,
 239 and P is the total number of instances according to the batch size used to train the
 240 CAE.

241 The MSE calculates the average of the squares of the discrepancies between the
 242 model’s predictions (\hat{X}_i) and the actual values (X_i) for a dataset composed of P
 243 examples. Its primary purpose is to minimize this value, aiming to make the outputs
 244 (\hat{X}_i) as similar as possible to the original inputs (X_i). This reflects the central goal
 245 of the CAE in learning a more concise and informative representation of the inputs,
 246 where a smaller value indicates a more precise reconstruction of the inputs.

247 However, the MSE term alone cannot guarantee diversity between the generated
 248 representations inside the ensemble. Consequently, a custom penalty term is added to
 249 the loss function after generating the first representation to ensure diversity between
 250 representations. Since we want to maximize the difference between the latent represen-
 251 tations of different CAEs, the penalty term is added to the loss function as a negative
 252 term. Equation 2 represents the proposed custom loss function:

$$\mathcal{L}_{\text{MSEP}} = \frac{1}{P} \sum_{i=1}^P (X_i - \hat{X}_i)^2 - \beta \frac{1}{T} \sum_{k=1}^T (R_k^{\text{last}} - R_k^{\text{curr}})^2 \quad (2)$$

253 where β defines the penalty provided to the similarity between the output of the
 254 R_k^{last} and R_k^{curr} , i.e., the latent vectors generated by the current and previous CAE,
 255 which is averaged across a validation set comprised of T instances. The dimensionality
 256 of the latent vectors R_k^{last} and R_k^{curr} may differ due to the differences in the CAE
 257 architectures, accordingly with the strategy used to generate diversity in the pool.
 258 This is the case when varying the dimensionality of the latent representation is the
 259 used strategy. Thus, we applied the principal component analysis (PCA) to make both
 260 latent vectors have the same dimensionality to allow computing the penalty term in
 261 Equation 2. For this purpose, we consider the size of the smallest latent vector as the
 262 number of selected PCA components.

263 The greater the similarity between the latent vectors, the greater the penalty
 264 applied. The idea is to increase the diversity between the latent representations of a
 265 previous and a current CAE at each iteration. The need to adjust the β parameter
 266 emerges from the search for a crucial trade-off between the representations generated
 267 and the desired diversity. The sweet spot is where the quality of the representations
 268 produced is maximized while encouraging the desired diversity.

269 3.3 Representation Ensemble Learning (REL) Method

270 This section introduces the proposed REL method, which is described here by three
 271 algorithms. Algorithm 1 presents the primary function responsible for generating a
 272 pool of unsupervised representations. The input parameters of the REL function play
 273 a pivotal role in configuring and executing the proposed method. In particular, we
 274 have ψ representing the number of representations to be generated, `str` defining the
 275 strategy to create the pool (S, L, A, SL, SA, LA, or SLA), X_t^u as the auxiliary
 276 unsupervised training dataset, X_v^u as the additional unsupervised validation dataset,
 277 both used for training the created autoencoders. The parameter β is related to the
 278 proposed custom loss function, defined in Equation 2 and used in Algorithm 2. It
 279 is a constant value that directly influences the magnitude and direction of the loss
 280 term in our function with the idea of penalizing the generation of similar CAEs.
 281 Finally, the structure `params` contains the configuration parameters for a default CAE
 282 architecture, including the number of layers, filter sizes, activation functions, and other
 283 fundamental hyperparameters. The default CAE is used if `params` is NULL.

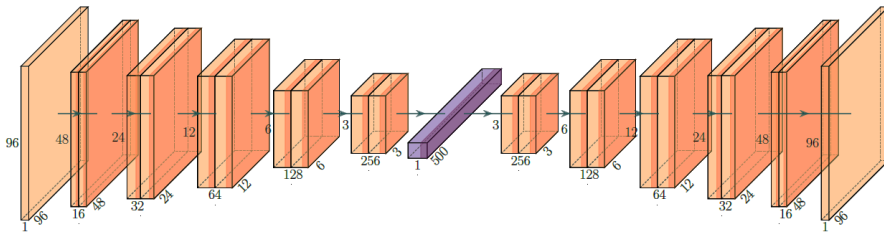


Fig. 3 Default CAE Architecture.

284 Fig. 3 shows the default CAE architecture used. It has input and output layers of
 285 dimensionality $96 \times 96 \times 1$, a depth of 5, and a latent vector size equal to 500. However,
 286 depending on the strategy used to generate the pool, the latent vector dimension and
 287 the CAE depth may change at each iteration of the REL function, according to lines
 288 9-13 and 14-20 of the Algorithm 1, respectively. All the parameter values of the default
 289 CAE are shown in Table 1. Such parameters were defined empirically based on prelim-
 290 inary experiments and iterative adjustments performed throughout the development
 291 of the proposed model. We have experimentally defined these CAE parameter values
 292 with FER performance in mind. It is worth noticing that the CAE can be customized
 293 to suit specific application requirements. Finally, the output parameter ϱ is a pool of
 294 feature extractors trained on an unsupervised approach.

295 At each iteration of the Algorithm 1, we create an encoder, a latent vector, and a
 296 decoder to compose a new CAE (see lines 21-28), which is compiled at line 29 consid-
 297 ering the provided `Optimizer` and `CustomLoss` parameters. Following that, on lines
 298 30-34, the created CAE is trained. As one can see in line 33, from the second iterat-
 299 ion, the output prediction (R_k^{last}) of the previous CAE estimated on the validation
 300 dataset (X_v^u) at line 35 is available to compute the customized loss function correctly.
 301 In line 36, we append the encoder in the pool as a new member, a feature extractor.

Algorithm 1 Representation Ensemble Learning (REL)

Input: ψ , as the number of representations,
 X_t^u , as the auxiliary train dataset,
 X_v^u , as the auxiliary validation dataset,
params, as the default CAE parameters (Table 1),
 β , as a constant for loss penalization term,
str, as the strategy to vary CAE parameters

Output: ϱ , as the pool of diverse representations

```
1: function REL( $\psi$ , str,  $X_t^u$ ,  $X_v^u$ ,  $\beta$ , params)
2:   max_depth  $\leftarrow$  params.CAE_Depth
3:   latent_size  $\leftarrow$  params.Latent_Size
4:   seed_number  $\leftarrow$  params.Seed
5:   for  $k$  in [1.. $\psi$ ] do
6:     if (str in [S, SL, SA, SLA]) then
7:       seed_number  $\leftarrow$  random(range(0 : params.Max_random_seed))
8:     end if
9:     if str in [L, SL, LA, SLA] then
10:      ini  $\leftarrow$  params.Min_latent_size
11:      end  $\leftarrow$  params.Max_latent_size
12:      latent_size  $\leftarrow$  random(range(ini : end))
13:    end if
14:    if (str in [A, SA, LA, SLA]) then
15:      if (max_depth = 1) then
16:        max_depth  $\leftarrow$  params.CAE_Depth
17:      end if
18:    else
19:      max_depth  $\leftarrow$  max_depth - 1
20:    end if
21:    for  $i$  in [1..max_depth] do
22:      encoder.create_convolutional_layer(params)
23:    end for
24:    encoder.create_flatten_layer(params)
25:    for  $i$  in [1..max_depth] do
26:      decoder.create_transpose_layer(params)
27:    end for
28:    CAE.(encoder, decoder) ▷ CAE model constructor
29:    CAE.compile(params.Optimizer, params.Custom_Loss)
30:    if  $k = 1$  then
31:      CAE.fit(params.Epoch, params.Batch_Size,  $X_t^u$ ,  $X_v^u$ , null)
32:    else
33:      CAE.fit(params.Epoch, params.Batch_Size,  $X_t^u$ ,  $X_v^u$ ,  $R_k^{\text{last}}$ )
34:    end if
35:     $R_k^{\text{last}}$   $\leftarrow$  CAE.encoder.predict( $X_v^u$ ) ▷ Prediction on validation set
36:     $\varrho$ .append(CAE.encoder)
37:    delete CAE
38:  end for
39:  return  $\varrho$ 
40: end function
```

Table 1 Default parameter values for the CAE model.

Parameter	Value(s)
Stride	2
Activation	RELU
Kernel_Size	3
Filters	[16, 32, 64, 128]
Output_Activation	Linear
Input_Size_Data	[96, 96, 1]
Latent_Size	500
Seed	42
CAE_Depth	5
Custom_Loss	Custom Loss (Equation 2)
Optimizer	SGD
Epoch	20
Batch_Size	60
Padding	Same
Max_random_seed	10000
Min_latent_size	150
Max_latent_size	2500

302 We detail the `Custom_Loss` function in the Algorithm 2. As one can see in line 2 of
 303 this algorithm, in the first iteration, REL uses a regular loss denoted on the Equation
 304 1. However, for the second representation and so on, the cost function used is the
 305 customized one, indicated in Equation 2. When training a new representation, we try
 306 to diversify regarding the one created in the last REL iteration.

307 In line 5 of Algorithm 2, the encoder corresponding to the current representation
 308 extracts the features from the validation set X_v^u . This is done using the predict func-
 309 tion, and the result obtained is stored in R_k^{curr} . If we opt for the latent vector strategy
 310 to introduce diversity in the representations, it is essential to highlight that the dimen-
 311 sionality of the intermediate layer will vary between them. Therefore, for us to correctly
 312 calculate the difference between R_k^{curr} and R_k^{last} , it is essential to equalize their dimen-
 313 sions. Thus, in line 6, the smallest value between these dimensions is obtained and used
 314 as the number of components to be considered in the PCA, as indicated in lines 7 and
 315 8. This process of equalizing the dimensions is crucial for guaranteeing the consistency
 316 and comparability of representations throughout the algorithm.

317 In line 9 of Algorithm 2, we calculate the sum of the quadratic differences between
 318 the last generated and current representations. The result of this calculation is then
 319 used, on line 10, as part of the penalty term, multiplying the value of β . This process
 320 is essential to adjust the impact of the penalty on the cost function and ensure that
 321 the difference between representations is considered according to the regularization
 322 parameter β . Such a regularization parameter reflects our concern about generating
 323 similar representations.

324 After completing the training of all representations as outlined in the Algorithm
 325 1, we move on to the next step, called “Feature Building” (Step 2 of Fig. 1), which
 326 extracts feature vectors from the labeled data X^l using the weights of each previously
 327 trained `CAE.encoder` available at the pool ϱ . We detail this process in the Algorithm
 328 3, which has the labeled target dataset X^l and the pool of feature extractors ϱ as

Algorithm 2 Custom Loss Function - Eq. 2

Input: β , as the constant value for the loss function,
 X_t^u , as the auxiliary train dataset,
 X_v^u , as the auxiliary validation dataset,
 R_k^{last} , as the output of the latent layer of the last CAE,
 X_i , as the true pattern,
 \hat{X}_i , as the output pattern

Output: loss_value, as the computed loss value

```
1: function CUSTOM_LOSS( $\beta$ ,  $X_t^u$ ,  $X_v^u$ ,  $R_k^{last}$ )
2:   if ( $R_k^{last} = \text{NULL}$ ) then
3:     loss_value  $\leftarrow \frac{1}{P} \sum_{i=1}^P (X_i - \hat{X}_i)^2$  (Equation 1)
4:   else
5:      $R_k^{curr} \leftarrow \text{CAE.encoder.predict}(X_v^u)$ 
6:     n_components  $\leftarrow \min(\text{length}(R_k^{last}), \text{length}(R_k^{curr}))$ 
7:     PCA_curr  $\leftarrow \text{PCA}(\text{n\_components}).\text{fit\_transform}(R_k^{last})$ 
8:     PCA_last  $\leftarrow \text{PCA}(\text{n\_components}).\text{fit\_transform}(R_k^{curr})$ 
9:     dif  $\leftarrow \sum_{k=1}^T (\text{PCA\_last} - \text{PCA\_curr})^2$ 
10:    penalty_term  $\leftarrow \beta \times \text{dif}$ 
11:    loss_value  $\leftarrow \frac{1}{P} \sum_{i=1}^P (X_i - \hat{X}_i)^2 - \text{penalty\_term}$  (Equation 2)
12:  end if
13:  return loss_value
14: end function
```

329 input. The output is a pool of feature sets A^l considered a novel way to represent
330 X^l . Such a new representation of the target dataset can subsequently be employed in
331 constructing supervised classifiers.

Algorithm 3 Feature Building Step 1

Input: X^l as the target dataset,
 ϱ as feature extractor pool

Output: A^l as a pool of new labeled features of the target dataset

```
1: function FEATURE_EXTRACTION( $X^l$ ,  $\varrho$ )
2:    $A^l \leftarrow []$ 
3:   for  $fe$  in [ $\varrho$ ] do
4:     feature_set  $\leftarrow fe.\text{predict}(X^l)$ 
5:      $A^l.\text{append}(\text{feature\_set})$ 
6:   end for
7:   return  $A^l$ 
8: end function
```

332 The supervised model training phase represents the final step of the self-taught
333 learning process. In this step, we evaluate the performance of a diverse set of classifiers,
334 including support vector machines (SVM), random forest (RF), bagging (BG) with

335 decision tree as a base classifier, and Knora union-based dynamic selection (KnoraU) in
 336 the context of decision trees (DT) and RF ensembles. These classifiers were examined
 337 to determine their effectiveness in leveraging the previously generated representations.
 338 The evaluated fusion strategies include sum, product, and stacking. **Table 2 presents**
 339 **the parameters empirically defined for each technique evaluated in the last step of the**
 340 **proposed method. The parameters not present in that table use the default values**
 341 **available in the Scikit-learn framework [34]. We conducted all experiments considering**
 342 **the same configuration, ensuring the reproducibility of the observed results.**

Table 2 Parameter settings of algorithms used in Step 3. Classifiers: single (SVM), ensembles (BG: bagging with decision tree as a base classifier, RF: random forest), KnoraU (DT: pool of trees; RF: random forest).

Algorithm	Parameters	
SVM	Kernel Penalty Parameter (C) Class Weight Probability	Linear 1e-6 Balanced True
BG with DT	Max Depth Tree Max Features Number of Base Estimators % of Training Samples	10 sqrt 100 1.0
RF	Max Depth Number of Trees Oob_score	10 100 True
KnoraU	Number of neighbors pool_classifiers	7 Bagging DT RF
Stacking	Meta classifier Solver Penalty parameter C	Logistic Regression lbfgs 0.1

343 4 Experiments and Discussion

344 This section describes the experiments to evaluate the proposed REL method,
 345 answering our two research questions.

346 4.1 Experimental Protocol

347 The target datasets are the Japanese Female Facial Expression (JAFFE) [35] and the
 348 Extended Cohn-Kanade (CK+) [36]. The JAFFE dataset is a laboratory-controlled
 349 image dataset with 213 images of 10 subjects (Japanese female models), with six basic
 350 facial expressions (happiness, anger, disgust, fear, sadness, and surprise) and a neutral
 351 one. In the JAFFE dataset, all ten subjects have one or more images for each class.
 352 The CK+ dataset is a laboratory-controlled dataset widely used to evaluate FER
 353 systems. It contains 523 sequences from 123 subjects, and 327 images of 118 subjects

354 are labeled with seven basic facial expressions (anger, contempt, disgust, fear, joy,
 355 sadness, and surprise). In the CK+ dataset, not all subjects have images in all classes.

356 Before extracting features from face images, we pre-processed them, cropping only
 357 the region of the region of interest (face) and selecting the reference points. Fig. 4 illus-
 358 trates the pre-processing applied to a sample image from JAFFE and CK+ datasets.
 359 To detect and crop the face area of the image, we use the Viola-Jones face detection
 360 method [37]. The reference points on the cropped image align the images so they are
 361 in the same position.

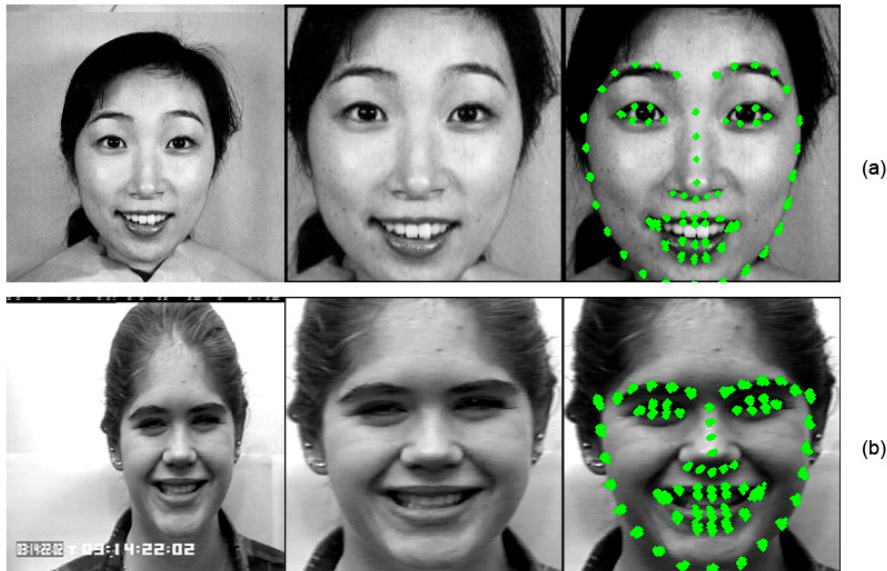


Fig. 4 Example illustrating the original image (left), face detection (middle), and landmark extraction (right) for (a) the JAFFE dataset and (b) the CK+ dataset.

362 As mentioned before, an unlabeled auxiliary dataset is necessary. The following
 363 auxiliary datasets were considered in our experiments:

- 364 • Kyoto Natural Images [38]: The dataset contains 62 natural images of resolution
 365 256×200 pixels, used by several other works that apply the concepts of STL, thanks
 366 to its large variability. The images from our experimental protocol belong to a
 367 domain (FER) far from that dataset.
- 368 • LabelMe [39]: The dataset consists of 50,000 images (40,000 for training and 10,000
 369 for testing), and each image is 256×256 pixels in size. The images belong to one
 370 of the 12 object classes: person, car, building, window, tree, sign, door, bookshelf,
 371 chair, table, keyboard, and head. For this paper, we used only the testing base.
- 372 • Labelled Face in the Wild (LFW) [40]: is a publicly available dataset of face pho-
 373 tographs used for verification, also known as peer matching. The dataset has 13,233
 374 images of faces collected on the web from 5,749 people.

375 The first two auxiliary datasets have images that do not belong to the FER domain
 376 but differ in size (62 images for Kyoto and 10,000 for LabelMe). The rationale is to
 377 show whether the size of the auxiliary dataset may impact the proposed method. On
 378 the contrary, the third dataset comprises faces closer to the target domain. Fig. 5
 379 shows four image samples of each auxiliary dataset used in our experiments.

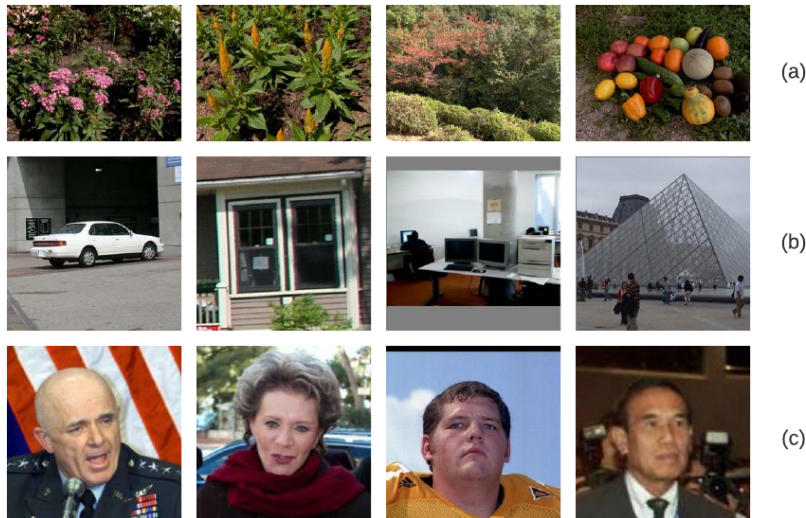


Fig. 5 Samples from the datasets (a) Kyoto Natural Images, (b) LabelMe, and (c) Labelled Face in The Wild.

380 The experimental protocol employs a LOSO cross-validation in the last step of the
 381 proposed method. In this cross-validation approach, the dataset is partitioned so that
 382 no subject in the test set appears in the training set. Suppose we have a dataset with
 383 N subjects; therefore, for each i -th fold, one subject will be designated as a test and
 384 the others for training. In such a protocol, when using the KnoraU dynamic ensemble
 385 selection, a subset of data (validation set) is extracted from the training set. Such a
 386 validation set is necessary since KnoraU needs to compute the competence of each
 387 classifier.

388 4.2 Results on the FER problem

389 For each target dataset (JAFPE and CK+), we performed experiments considering the
 390 three auxiliary datasets (Kyoto, LabelMe, and LFW) and representation ensembles of
 391 different sizes (5, 10, 15, 20, 30, 40, 50, 70, 100, and 150). The strategies S, A, L, SA, SL,
 392 LA, and SLA were tested in each experiment. It is crucial to notice that strategy A was
 393 used uniquely in the investigation, where we generated just five representations because
 394 it explores the depth of the CAE architecture. Thus, using a depth greater than

395 five layers makes the experiment almost impractical on the computational resources
396 available for running the experiments.

397 At the last step of our STL-based approach, we used the following inducers: SVM,
398 RF, BG, and dynamic ensemble selection based on KnoraU executed in a pool of DTs
399 and RF. For each ensemble of representations, we present the result of the best rep-
400 resentation and their fusion (sum, product, and stacking). Furthermore, the standard
401 deviation of all models inside the ensemble is calculated to provide an idea of the
402 variability achieved in the generated ensemble.

403 4.2.1 JAFFE as Target Dataset

404 Tables 3, 4, and 5 present the results for the JAFFE dataset when CAE is trained with
405 Kyoto, LabelMe, and LFW datasets respectively. The best outcome for the JAFFE
406 datasets was achieved with an ensemble of 50 representations trained on the Kyoto
407 dataset, varying the latent vector dimensionality (strategy L), and using KnoraU
408 dynamic ensemble selection executed on an RF at each representation. Thus, by com-
409 bining the results of each representation using stacking, we achieved 66.66% accuracy.
410 The standard deviation inside the corresponding pool of representations was (+/-)
411 3.33 percentage points. It is also important to note that the best single representation
412 using KnoraU provided 42.00% accuracy.

413 When using the LabelMe dataset (Table 4), the best result was achieved with an
414 ensemble composed of 50 representations varying the latent vector size, architecture,
415 and seeds (strategy SLA). The last step (supervised) uses KnoraU dynamic ensemble
416 selection executed on an RF at each representation. The best single representation
417 provided 44.67% of accuracy, while by combining the results of each representation
418 using stacking, we achieved 65.25%.

419 When using the LFW dataset (Table 5), the best result was achieved with an
420 ensemble composed of ten representations. The best accuracy (65.25%) was observed
421 by varying the seeds and the CAE architecture (strategy SA). The final inducer is
422 an ensemble of SVMs combined via product rule. The best single SVM inside the
423 ensemble achieved 61.39% of accuracy.

424 Tables 3, 4 and 5 also present the accuracy standard deviation among generated
425 representations (between brackets beside the best single result) that varies from 2.50 to
426 6.89 percentage points (4.02 in average), what was expected in our method. **The result
427 analysis indicates that after 50 representations, there is no significant improvement
428 in performance. Such an observation shows that, as in an ensemble of classifiers, we
429 must empirically define the number of representations in our pool for each problem.
430 In other words, a larger ensemble does not always mean a better performance. For
431 instance, Fig. 6 shows that considering the results produced with the Kyoto dataset,
432 an ensemble with 50 members may provide the same result as an ensemble with 150
433 members, substantially reducing complexity.**

434 4.2.2 CK+ as Target Dataset

435 Tables 6, 7, and 8 present the results for the CK+ dataset when CAE is trained
436 with Kyoto, LabemMe, and LFW datasets, respectively. The best result for the CK+
437 dataset (92.66%) was achieved with an ensemble of 50 representations trained on the

Table 3 Accuracy results using the LOSO protocol with a pool of 50 representations generated by the proposed method. Kyoto and JAFFE are used as auxiliary and target datasets, respectively. FER models based on a pool of SVMs, RF, BG, and dynamic selection based on KnoraU in a collection of DTs and RF. Strategies S, L, and A represent different seeds, latent layer dimensions, and CAE architectures. For each ensemble, the single best classifier and the following fusion strategies: sum, product, and stacking. (*) is the pool accuracy standard deviation. The best results are in bold.

Strategy		Pool of SVMs	Ensembles		Dynamic Selection (KnoraU)	
			BG	RF	DT	RF
S	Best Single	58.50 (3.24)	55.31 (4.24)	40.76 (3.63)	56.67 (4.22)	45.18 (4.12)
	Sum	59.15	61.50	60.09	59.15	58.21
	Product	59.15	57.27	59.62	59.15	57.27
	Stacking	59.62	58.68	43.66	61.97	63.38
L	Best Single	60.42 (4.00)	51.37 (4.26)	44.74 (3.94)	51.89 (4.30)	42.00 (3.33)
	Sum	60.09	61.03	60.56	61.03	62.44
	Product	60.09	61.97	60.56	61.03	62.91
	Stacking	61.03	60.56	48.82	61.97	66.66
SA	Best Single	61.19 (3.19)	49.54 (4.05)	44.45 (3.60)	50.47 (4.14)	42.71 (3.57)
	Sum	61.97	61.32	62.91	59.62	61.50
	Product	61.97	60.56	64.78	60.09	61.97
	Stacking	60.56	59.62	53.05	60.56	63.84
SL	Best Single	61.47 (3.41)	52.63 (4.40)	43.50 (3.74)	51.72 (4.21)	44.38 (4.40)
	Sum	61.03	61.97	57.74	61.03	60.09
	Product	61.50	60.56	57.27	62.91	60.09
	Stacking	62.44	59.62	43.66	62.91	63.84
LA	Best Single	60.61 (3.38)	50.05 (3.38)	47.23 (4.45)	52.88 (3.52)	48.14 (4.90)
	Sum	60.56	61.03	61.97	61.03	62.91
	Product	60.56	60.09	62.44	61.50	62.91
	Stacking	60.03	61.97	51.17	61.50	61.50
SLA	Best Single	61.35 (3.37)	51.62 (4.36)	44.46 (4.36)	52.94 (4.11)	47.41 (4.20)
	Sum	61.50	62.91	59.15	63.38	59.15
	Product	61.50	61.97	59.62	61.97	59.62
	Stacking	61.50	58.68	48.35	60.09	62.44

438 Kyoto dataset, varying the seeds and latent vector dimensionality (strategy SL) and
439 using the fusion based on the stacking of SVMs. It is also essential to note that the
440 best result based on a single SVM classifier in this experiment was 91.29%.

441 When using the LabelMe dataset (Table 7), the best result of 91.74% was achieved
442 with an ensemble composed of 30 representations and by varying the latent vector
443 dimensionality (strategy L), using an ensemble of SVMs. The best result when using
444 the LFW dataset as an auxiliary dataset (90.51%) was observed using the SL strategy
445 and an ensemble of SVMs, as shown in Table 8.

446 Unlike the JAFFE dataset, ensemble methods (BG, RF) and dynamic selection
447 (KnoraU) were less effective on the CK+ dataset. One possible explanation is the
448 relatively small standard deviation in accuracy among the generated representations
449 on CK+ compared to JAFFE. The standard deviation for CK+ ranged from 0.7 to
450 5.4 percentage points (average of 2.3), suggesting a lack of significant variation in the
451 quality of the representations.

452 As observed in the previous experiments on the JAFFE dataset, the behavior is
453 similar here, confirming that we must define the best ensemble size empirically. Fig. 7
454 shows that the best ensemble size was 50 when considering Kyoto and 70 when using

Table 4 Accuracy results using the LOSO protocol with a pool of 50 representations generated by the proposed method. LabelMe and JAFFE are used as auxiliary and target datasets, respectively. FER models based on a pool of SVMs, RF, BG, and dynamic selection based on KnoraU in a collection of DTs and RF. Strategies S, L, and A represent different seeds, latent layer dimensions, and CAE architectures. For each ensemble, the single best classifier and the following fusion strategies: sum, product, and stacking. (*) is the pool accuracy standard deviation. The best results are in bold.

Strategy		Pool of SVMs	Ensembles		Dinamic Selection (KnoraU)	
			BG	RF	DT	RF
S	Best Single	57.08 (3.32)	50.53 (4.19)	44.30 (5.23)	50.79 (4.52)	42.44 (4.77)
	Sum	55.86	61.97	59.15	61.03	61.03
	Product	57.27	61.03	59.62	62.44	60.56
	Stacking	56.33	59.15	50.70	62.44	62.44
L	Best Single	58.96 (5.56)	48.66 (4.36)	42.62 (3.84)	49.58 (4.30)	43.33 (3.59)
	Sum	59.15	56.33	58.21	56.33	57.27
	Product	58.68	56.33	58.21	55.86	56.80
	Stacking	58.68	55.86	42.25	58.21	60.56
SA	Best Single	61.36 (3.20)	49.58 (3.96)	43.97 (4.72)	49.06 (3.94)	44.30 (4.24)
	Sum	60.56	60.56	62.91	60.09	63.84
	Product	61.50	58.68	62.91	60.56	63.38
	Stacking	58.68	62.44	53.52	59.62	64.78
SL	Best Single	59.44 (2.98)	46.35 (3.11)	44.18 (3.90)	47.86 (3.12)	44.31 (3.54)
	Sum	58.21	61.50	60.09	59.62	59.15
	Product	57.74	55.86	60.56	59.62	59.15
	Stacking	59.62	59.15	46.94	63.38	62.44
LA	Best Single	60.40 (3.51)	50.82 (5.00)	45.58 (4.65)	48.72 (4.96)	43.41 (4.61)
	Sum	59.62	57.27	59.15	59.62	58.68
	Product	60.09	58.68	60.09	58.68	58.21
	Stacking	60.56	61.00	47.41	61.03	62.44
SLA	Best Single	60.78 (3.89)	50.15 (3.98)	44.12 (3.84)	52.52 (4.05)	44.67 (4.14)
	Sum	60.56	59.62	60.09	60.56	60.56
	Product	61.03	61.03	59.62	60.56	60.56
	Stacking	59.62	61.03	50.23	61.97	65.25

455 LabelMe and LFW datasets. It points out that the performance of the generated
456 ensemble is relative to the capacity to cover the problem space with complementary
457 models, which does not depend on the pool size but on the algorithm’s ability to
458 generate diverse representations.

459 4.3 Ablation Study

460 After generating the first representation, we consider a penalty term in the cost func-
461 tion to ensure that subsequent representations are diverse, as denoted in Equation 2.
462 The β parameter is a user-defined constant that plays an essential role in this sce-
463 nario, as it is used to control the penalty applied for generating diversity between the
464 representations. In other words, it regulates the importance of the penalty applied to
465 representations to encourage the exploration of different regions of the feature space.
466 The idea is to generate different autoencoders at each training step. We opted for a
467 uniform β in this work, which proved effective and kept the model simple and effi-
468 cient. However, different strategies can be investigated, for instance, varying the value

Table 5 Accuracy considering the LOSO protocol and a pool with 10 representations generated with the proposed method. LFW and JAFFE are used as auxiliary and target datasets, respectively. FER models based on a pool of SVMs, RF, BG, and dynamic selection based on KnoraU in a collection of DTs and RF. Strategies S, L, and A represent different seeds, latent layer dimensions, and CAE architectures. For each ensemble, the single best classifier and the following fusion strategies: sum, product, and stacking. (*) is the pool accuracy standard deviation. The best results are in bold.

Strategy		Pool of SVMs	Ensembles		Dynamic Selection (KnoraU)	
			BG	RF	DT	RF
S	Best Single	60.52 (3.65)	44.99 (3.23)	38.93 (2.65)	42.72 (3.08)	38.73 (2.67)
	Sum	58.68	53.99	50.70	52.58	51.64
	Product	60.09	53.52	51.64	53.52	51.17
	Stacking	60.56	55.39	44.13	51.17	53.52
L	Best Single	60.56 (3.82)	48.04 (3.84)	47.26 (5.11)	49.42 (3.96)	47.26 (4.83)
	Sum	61.03	58.21	56.80	58.21	61.97
	Product	60.56	60.56	57.74	58.68	61.03
	Stacking	59.15	56.80	48.82	56.33	60.09
SA	Best Single	61.39 (2.99)	46.33 (5.04)	39.72 (3.35)	50.34 (4.12)	42.87 (3.16)
	Sum	64.31	58.21	58.21	57.27	57.74
	Product	65.25	58.68	60.09	57.27	58.68
	Stacking	63.38	58.21	51.17	54.92	53.05
SL	Best Single	58.53 (4.01)	42.44 (3.78)	39.54 (3.18)	42.72 (3.64)	40.44 (3.51)
	Sum	60.09	55.86	57.74	52.58	55.86
	Product	59.62	57.74	56.80	52.58	55.39
	Stacking	59.62	54.92	46.94	49.29	55.86
LA	Best Single	59.17 (2.50)	50.83 (4.86)	43.97 (4.87)	51.21 (5.31)	45.38 (5.22)
	Sum	58.68	57.27	56.33	56.33	54.92
	Product	59.62	58.21	56.33	57.74	54.46
	Stacking	59.15	58.68	52.11	56.33	57.74
SLA	Best Single	58.92 (3.02)	54.24 (6.89)	42.54 (5.32)	51.88 (6.38)	44.51 (4.66)
	Sum	61.97	56.80	59.15	54.92	58.21
	Product	62.44	54.92	57.27	55.86	57.27
	Stacking	61.97	56.33	52.58	53.05	58.21

469 of β across the training steps. Increasing β value may be a strategy to create more
470 different autoencoders during training.

471 The best results for JAFFE and CK+ datasets are shown in Tables 3 and 6,
472 respectively. For the ablation study, we used the best performance settings. Thus, we
473 generated 50 representations with the unlabeled dataset Kyoto and applied different
474 values for the β . Table 9 shows the performance of the obtained model for each value
475 of β .

476 The best results on JAFFE and CK+ datasets were found using $\beta = 0.001$. When
477 we increase or decrease the value of β , we reduce model performance. The β value is
478 experimentally adjusted to find the right balance between the quality of the generated
479 representations and the desired diversity. If the β is too low, diversity might not be
480 encouraged enough, and representations might become too similar. On the other hand,
481 if the β is too high, the penalty can completely overwhelm the objective, resulting in
482 very diverse but potentially low-quality representations.

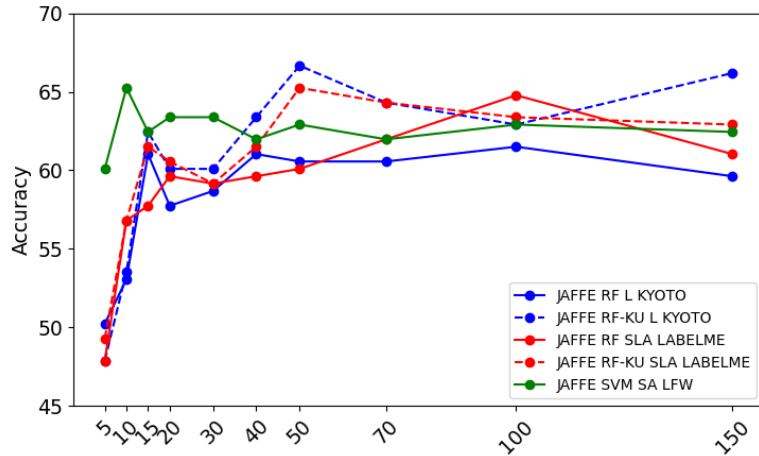


Fig. 6 The impact of varying the ensemble size in the best setup achieved for the JAFFE dataset when using Kyoto, LabelMe, and LFW datasets as auxiliary ones.

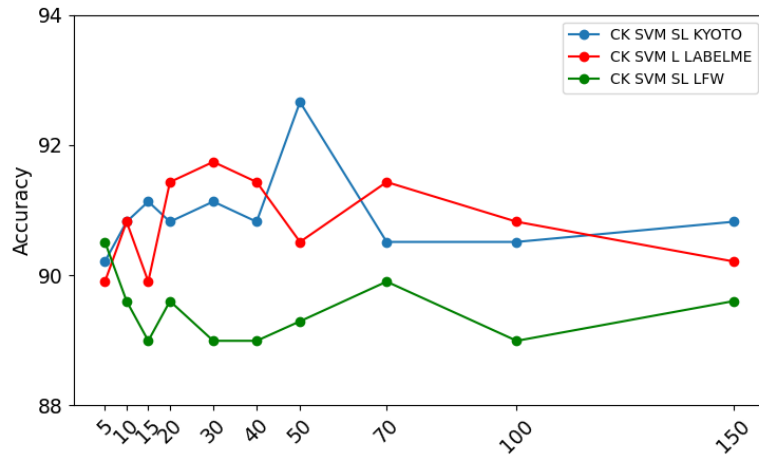


Fig. 7 The impact of varying the ensemble size in the best setup achieved for the CK+ dataset when using Kyoto, LabelMe, and LFW datasets as auxiliary ones.

483 **4.4 Comparison with CNN-based Approaches and Related**
 484 **Works**

485 As stated in Section 1, the performance of supervised approaches against STL is a
 486 matter of discussion once the former is widely used in the literature and generally
 487 achieves the best results. Thus, to answer our RQ2 correctly, we assess the traditional
 488 supervised learning strategy using the LOSO protocol presented in Section 4.1, in
 489 which the samples of test subjects are not present in the training step.

Table 6 Accuracy considering the LOSO protocol and a pool with 50 representations generated with the proposed method. Kyoto and CK+ are used as auxiliary and target datasets, respectively. FER models based on a pool of SVMs, RF, BG, and dynamic selection based on KnoraU in a pool of DTs and RF. Strategies S, L, and A represent different seeds, latent vector dimensions, and CAE architectures. For each ensemble, the single best classifier and the following fusion strategies: sum, product, and stacking. (*) is the pool accuracy standard deviation. The best results are in bold.

Strategy		Pool of SVMs	Ensembles		Dinamic Selection (KnoraU)	
			BG	RF	DT	RF
S	Single Best	88.48 (2.63)	71.75 (3.09)	65.09 (2.01)	71.90 (2.79)	66.62 (1.90)
	Sum	87.15	74.00	68.19	73.39	68.80
	Product	87.46	74.92	69.11	74.31	68.80
	Stacking	89.90	80.42	76.75	76.45	80.12
L	Single Best	89.87 (2.58)	72.45 (2.54)	69.78 (2.34)	71.10 (2.26)	69.13 (2.30)
	Sum	89.60	75.84	71.55	75.22	72.17
	Product	89.29	76.45	71.55	76.14	72.17
	Stacking	90.82	77.98	78.89	74.31	81.34
SA	Single Best	89.47 (1.80)	71.55 (2.47)	66.66 (1.90)	73.88 (2.60)	67.74 (1.81)
	Sum	89.90	74.92	70.64	75.22	72.17
	Product	89.29	77.37	70.64	75.84	72.17
	Stacking	90.21	80.12	75.84	75.84	78.89
SL	Single Best	91.29 (2.86)	72.57 (2.72)	68.31 (2.14)	72.14 (2.82)	67.37 (1.85)
	Sum	89.60	74.92	70.64	74.00	71.55
	Product	89.60	75.84	70.64	74.00	71.55
	Stacking	92.66	77.67	77.06	74.31	78.89
LA	Single Best	87.18 (2.24)	72.06 (2.72)	68.88 (2.02)	73.33 (2.85)	69.51 (2.24)
	Sum	86.85	77.37	71.55	77.06	72.47
	Product	87.46	77.37	72.17	77.06	72.47
	Stacking	88.07	79.81	76.14	78.28	79.81
SLA	Single Best	87.00 (3.05)	74.23 (3.16)	68.14 (2.20)	73.29 (3.04)	68.58 (2.28)
	Sum	87.76	76.45	71.86	75.53	71.86
	Product	87.76	77.06	71.86	76.14	72.17
	Stacking	88.07	79.51	76.45	76.45	79.81

490 In the state-of-the-art, there are a plethora of ready-to-use CNN architectures. To
491 avoid a biased comparison, the rationale is to select well-known networks that provided
492 a breakthrough in their architectures and achieved competitive performances on the
493 ImageNet Challenge [41].

494 The first successful approach was made by Krizhevsky et al. [42], implementing
495 a convolutional architecture named AlexNet composed of five layers. Next, Simonyan
496 and Zisserman [43](VGG) proposed a more profound architecture. Contrary to the
497 rationale of stacking layers that increase the depth of the network, Szegedy et al. [44]
498 proposed the inception modules composed of a pattern of convolutional layers, pooling,
499 and feature concatenation. The final architecture is a stack of Inception modules that
500 provide a more comprehensive network and enhance the feature representation. So far,
501 most of the proposed architectures have resorted to the fact that high performances
502 should be achieved as the network grows. However, the vanishing of the gradient has
503 become an issue. To minimize the overfitting, skip connections between adjacent layers
504 were proposed in ResNet [45] and improved in DenseNet [46], which extends the skip
505 connections between all network layers. Last, EfficientNet [47] proposed a technique
506 to provide model scaling to improve the efficiency of the network.

Table 7 Accuracy considering the LOSO protocol and a pool with 30 representations generated with the proposed method. LabelMe and CK+ datasets are used as auxiliary and target datasets, respectively. FER models based on a pool of SVMs, RF, BG, and dynamic selection based on KnoraU in a collection of DTs and RF. Strategies S, L, and A represent different seeds, latent vector dimensions, and CAE architectures. For each ensemble, the single best classifier and the following fusion strategies: sum, product, and stacking. (*) is the pool accuracy standard deviation. The best results are in bold.

Strategy		Pool of SVMs	Ensembles		Dinamic Selection (KnoraU)	
			BG	RF	DT	RF
S	Best Single	89.57 (4.67)	69.50 (2.23)	66.93 (2.05)	69.97 (2.22)	66.72 (2.23)
	Sum	87.46	73.70	70.03	73.08	70.03
	Product	87.76	74.00	70.03	73.70	70.03
	Stacking	89.90	78.89	77.37	75.22	79.20
L	Best Single	90.38 (2.25)	72.16 (2.27)	65.43 (1.47)	71.99 (1.99)	66.76 (1.78)
	Sum	89.90	76.45	69.74	75.53	70.64
	Product	89.90	75.84	69.72	76.45	70.64
	Stacking	91.74	81.34	78.59	74.31	77.98
SA	Best Single	88.87 (1.91)	69.19 (2.13)	67.09 (2.11)	70.55 (2.16)	68.37 (2.09)
	Sum	88.68	74.92	70.33	75.53	70.33
	Product	88.68	76.75	70.64	75.53	70.33
	Stacking	89.29	79.81	75.22	77.98	77.37
SL	Best Single	89.71 (2.60)	71.08 (2.45)	69.03 (2.16)	71.51 (2.17)	67.31 (1.63)
	Sum	88.68	74.61	70.03	74.31	70.33
	Product	88.99	74.61	70.33	74.92	70.03
	Stacking	90.51	79.81	75.22	75.53	77.37
LA	Best Single	89.39 (3.84)	72.37 (2.68)	67.92 (2.53)	71.86 (2.82)	69.49 (3.01)
	Sum	90.21	77.06	70.94	75.84	70.94
	Product	89.90	77.67	71.25	76.75	70.94
	Stacking	90.51	80.12	75.53	77.98	79.51
SLA	Best Single	90.45 (2.38)	71.17 (3.21)	66.38 (2.36)	73.13 (3.43)	67.93 (1.95)
	Sum	89.29	75.22	70.33	74.92	70.64
	Product	89.29	74.92	71.25	75.22	70.64
	Stacking	90.51	79.81	76.45	76.14	76.75

507 The architectures mentioned above and their variants in terms of the number of
508 layers were trained to JAFFE and CK+ datasets¹. To a fair comparison, we have also
509 included the encoder architecture of the CAE (Section 3.3). Except for the AlexNet and
510 the CAEs, we have used the pre-trained weights on the ImageNet dataset, considering
511 two different fine-tuning schemas. First, (1) all weights of convolutional layers are
512 frozen, and only fully connected layers are fine-tuned to the FER problem. The second
513 schema also fine-tunes the last convolutional layer (2).

514 The rationale for using pre-trained weights is based on the fact that FER datasets
515 do not have sufficient data to fine-tune deep models correctly. Besides, this strategy is
516 a common practice in state-of-the-art [48–51] and should provide meaningful insights
517 when comparing the TL with the STL strategy.

518 Table 10 shows the performance of all experiments. The columns denoted with (1)
519 and (2) represent the results using the respective fine-tuning schema for FER datasets.

¹All trained models are available for research purposes at [the hyperlink will be inserted in the final version]

Table 8 Accuracy considering the LOSO protocol and a pool with 5 representations generated with the proposed method. LFW and CK+ are auxiliary and target datasets, respectively. FER models based on a pool of SVMs, RF, BG, and dynamic selection based on KnoraU in a collection of DTs and RF. Strategies S, L, and A represent different seeds, latent vectors, and CAE architectures. For each ensemble, the single best classifier and the following fusion strategies: sum, product, and stacking. (*) is the pool accuracy standard deviation. The best results are in bold.

Strategy		Pool of SVMs	Ensembles		Dinamic Selection (KnoraU)	
			BG	RF	DT	RF
S	Best Single	83.24 (1.69)	66.75 (1.68)	63.22 (1.28)	66.96 (2.04)	64.77 (1.17)
	Sum	85.93	73.08	68.80	71.86	69.41
	Product	85.93	73.39	69.11	73.08	69.72
	Stacking	87.15	74.92	73.70	72.47	72.47
A	Best Single	84.91 (0.74)	66.82 (1.04)	65.38 (1.42)	67.20 (0.96)	66.37 (1.80)
	Sum	87.46	74.00	67.88	73.70	70.03
	Product	88.07	74.92	68.19	74.00	70.33
	Stacking	87.76	76.45	73.70	73.39	74.31
L	Best Single	84.27 (3.47)	68.62 (3.21)	64.26 (2.50)	69.15 (3.20)	64.50 (2.25)
	Sum	87.76	75.35	68.80	75.53	68.50
	Product	87.46	77.06	68.19	74.92	68.80
	Stacking	87.46	78.28	74.92	76.14	73.39
SA	Best Single	88.16 (3.35)	71.96 (3.05)	66.77 (1.56)	72.45 (3.01)	67.07 (2.36)
	Sum	87.76	72.17	70.64	71.55	71.55
	Product	87.76	73.08	70.64	71.55	71.55
	Stacking	88.99	74.61	75.84	71.55	74.31
SL	Best Single	85.69 (0.99)	68.75 (2.29)	66.96 (2.46)	68.77 (1.88)	68.26 (2.22)
	Sum	90.21	75.53	70.33	75.84	70.64
	Product	89.60	75.22	70.33	75.53	70.94
	Stacking	90.51	78.89	77.67	75.22	77.37
LA	Best Single	83.92 (1.42)	71.56 (2.68)	67.17 (1.30)	69.83 (2.02)	68.24 (1.72)
	Sum	87.15	72.47	70.94	73.08	70.33
	Product	86.85	73.39	70.33	72.78	70.33
	Stacking	87.76	76.75	77.06	71.86	75.22
SLA	Best Single	86.70 (5.41)	71.14 (4.09)	66.97 (2.69)	70.67 (3.90)	65.14 (1.86)
	Sum	86.54	74.92	69.41	73.08	69.41
	Product	85.62	75.22	69.72	74.61	69.41
	Stacking	86.85	74.61	73.70	73.70	74.00

Table 9 Ablation experiment of the β parameter considering the best setup of the proposed method for JAFFE and CK+ datasets.

	β	Accuracy (%)
JAFFE	0.0001	60.56
	0.001	66.66
	0.01	62.91
	0.1	60.56
CK	0.0001	91.13
	0.001	92.66
	0.01	89.60
	0.1	90.82

Table 10 CNN architectures trained with the LOSO protocol and two fine-tuning strategies (1, 2).

Architecture	CK(1)	CK(2)	JAFFE(1)	JAFFE(2)
VGG_16	78.9	46.8	40.6	14.1
VGG_19	81.4	61.1	44.5	14.5
InceptionV3	75.7	84.2	42.2	36.9
ResNet50	84.3	81.6	56.7	47.9
ResNet101	83.5	81.3	43.4	51.2
InceptionResnetV2	78.5	84.6	36.6	44.2
EfficientNetB0	80.6	86.3	38.3	49.0
EfficientNetB7	75.2	83.5	33.6	52.6
DenseNet121	81.1	81.1	54.9	45.4
DenseNet201	85.8	82.0	54.5	50.7
Encoder (CAE)*		69.3		39.3
Alexnet*		70.6		42.9
Ours		92.66		66.66

(*)The model was trained from scratch (no pre-trained weights)

520 Finally, the models marked with an asterisk represent those trained from scratch once
521 they had fewer layers to fine-tune.

522 The results allow us to draw some conclusions. The proposed REL surpassed pre-
523 trained CNN models. The best CNN-based solution for the JAFFE dataset achieved
524 56.7% using ResNet50, while the proposed STL-based solution achieved 66.66% on the
525 same dataset. A similar behavior was observed for the CK+ dataset. The best CNN
526 solution achieved 86.3%, while the proposed STL-based method achieved an accuracy
527 of 92.66%. It means that independently of the fine-tuning schema, the STL can surpass
528 the TL by a fair margin, providing a better representation even when learning features
529 from a different context.

530 The trade-off between classification time and accuracy is discussed when com-
531 paring STL and traditional supervised approaches. To clarify this, we evaluated the
532 classification time for each test set using the LOSO protocol. The results are depicted
533 in Figs. 8 and 9, for JAFFE and CK+ datasets, respectively. In such an analysis,
534 we compare the best ensemble generated by the REL algorithm against traditional
535 supervised learning approaches, previously discussed in Table 10. They are compared
536 considering execution time during the test (classification time), accuracy, and model
537 size in megabytes.

538 The proposed approach generally exhibits a better trade-off regarding classification
539 time versus accuracy. A practical example can be seen in Fig. 8, where the RF-KnoraU
540 approach outperforms the well-known networks VGG16, Inception, and ResNet. In
541 Fig. 9, the result of the REL algorithm obtained using a pool of SVMs on the CK+
542 dataset outperforms CNN architectures in terms of accuracy and most of them in
543 terms of classification time.

544 Finally, we present a comparative study with the results of our proposed model,
545 which addresses the generation of a pool of unsupervised representations and its rela-
546 tionship with the most recent advances in facial emotion recognition, using the same

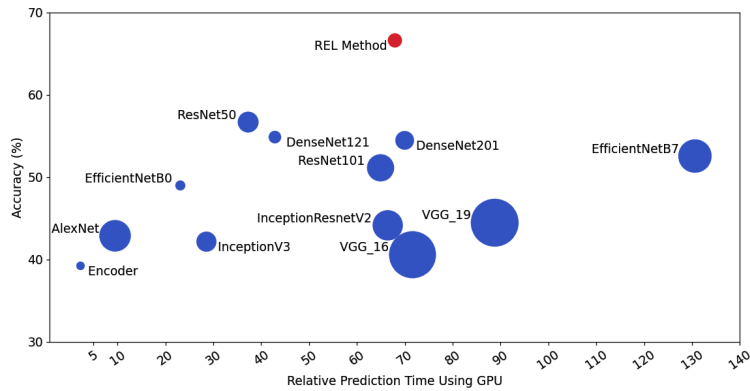


Fig. 8 Comparison of the proposed STL-based method (in red) with CNN architectures (in blue) in terms of accuracy (%), classification time (in seconds), and disk space (in megabytes) as the circle size. Results were computed using the JAFFE dataset and the LOSO protocol. The STL-based method considers the final classifier RF with KnoraU dynamic selection.

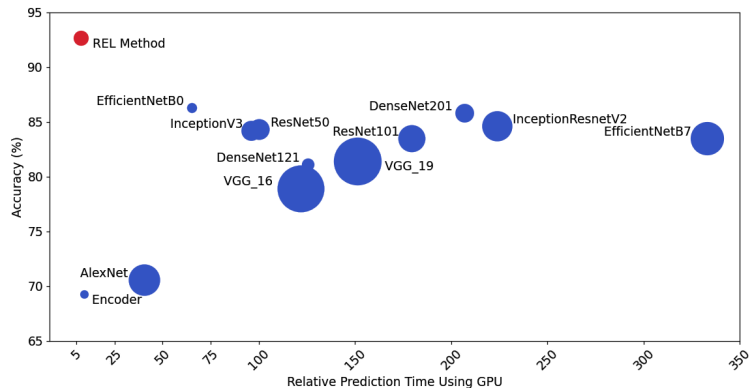


Fig. 9 Comparison of the proposed STL-based method (in red) with CNN architectures (in blue) in terms of accuracy (%), classification time (in seconds), and disk space (in megabytes) as the circle size. Results were computed using the CK+ dataset and the LOSO protocol. The STL-based method considers as the final classifier a pool of SVMs.

547 experimental protocol. This comparison aims to provide an expanded understanding
 548 of the contributions and improvements brought by our approach. Our effort focused
 549 on gathering the most significant possible number of studies that adhered to the same
 550 adopted evaluation protocol based on the LOSO protocol. Tables 11 and 12 present
 551 the highest accuracies observed in JAFFE and CK+ datasets, respectively, that have
 552 been documented in widely referenced scientific publications and conferences. The
 553 data revealed that the proposed method sets a new benchmark for facial expression
 554 recognition.

Table 11 Benchmark comparison on the JAFFE dataset using the LOSO protocol.

Reference	Method	Accuracy (%)	Classes	Feature Type
Kola & Samayamantula [52]	LGC-HD	60.70	6	Handcrafted
	LGC-HD	58.20	7	Handcrafted
Kartheek et al. [53]	RMP_Prime	61.64	6	Hand-crafted
Mandal et al. [54]	DRADAP	57.22	6	Handcrafted
	ARADAP	56.20	7	Handcrafted
	Our Method	66.66	7	Deep

Table 12 Benchmark comparison on the CK+ dataset using the LOSO protocol.

Reference	Method	Accuracy (%)	Classes	Feature Type
Du & Hu [55]	WPLBP	91.72	6	Handcrafted
	WPLBP	86.47	7	Handcrafted
Wu & Lin [56]	GM	86.83	7	Deep
	GM+AFM	87.78	7	Deep
	GM+ W-AFM	88.25	7	Deep
	GM + W-CRAFM	89.84	7	Deep
Lee et al. [57]	CER-ICV	92.34	7	Handcrafted
Kola & Samayamantula [52]	LGC-HD	72.80	6	Handcrafted
	LGC-HD	70.60	7	Handcrafted
	Our Method	92.66	7	Deep

555 Through these experiments, it becomes evident that generating a pool of unsu-
556 pervised representations characterized by their diversity plays a fundamental role in
557 improving the performance of an FER system.

558 4.5 Discussion

559 Concerning the auxiliary dataset, the experiments showed that the best results were
560 found when the ensemble of representations was learned using the small Kyoto dataset.
561 Moreover, such a dataset is far from the target domain, but it still provides better
562 results when compared to those observed using the LFW dataset, which is composed
563 of faces. **In summary, our results indicated that the variability of images within the
564 dataset is a more critical factor for performance than the total number of images or
565 the proximity to the target problem. This finding suggests that the right choice for the
566 auxiliary dataset must consider data diversity (image variability), which is essential
567 for learning robust representations.**

568 Based on the results observed in our experiments and analyses, we can posi-
569 tively answer our first research question, RQ1. The unsupervised representations with
570 diversity obtained through the REL algorithm significantly impacted the performance
571 of the proposed FER model. By incorporating diversity in an unsupervised man-
572 ner through different initialization strategies and a custom loss function, the REL
573 algorithm has demonstrated its effectiveness in producing a set of comprehensive
574 representations. These representations capture latent and complex features in facial
575 expressions, allowing the model to extract meaningful and distinctive features.

576 We obtained high accuracy in recognizing facial expressions by training the model
577 with the unsupervised representations generated by the REL. The proposed method
578 is up to 5.1% and 1.3% better than the single classifier for JAFFE and CK+ datasets,
579 respectively.

580 The diversity introduced by REL made it possible for the model to become more
581 adaptable to different variations in expressions, including nuances and subtle differ-
582 ences between categories of expressions. In addition, the generalizability of the model
583 has been greatly improved. The diversified representations allowed the model to learn
584 general characteristics that were not restricted to a specific set of training data, mak-
585 ing it more efficient in identifying facial expressions in new datasets. This reduces the
586 tendency to overfit by preventing the model from overspecializing on specific training
587 data.

588 Furthermore, including a pool of unsupervised representations generated by the
589 REL algorithm facilitated the knowledge transfer process. The pre-training of the
590 unsupervised representations allowed the FER model to initialize with more compre-
591 hensive knowledge and, later to be adjusted for the specific FER task. This improved
592 the overall performance of the model.

593 The answer to our second research question (RQ2) was only possible after the last
594 experiments in which the proposed STL-based approach was compared favorably with
595 the CNN-based FER solutions.

596 The comparative analysis between the proposed method and the CNN models
597 adjusted for JAFFE and CK+ datasets revealed considerable improvements in terms
598 of precision. Specifically, the proposed method outperformed the fitted CNN models
599 by up to 9.9% in the JAFFE dataset and by up to 6.3% in the CK+ dataset. These
600 substantial gains in accuracy are of great relevance and demonstrate the effectiveness
601 of the STL-based approach.

602 A favorable comparison with supervisory-based CNN solutions suggests that the
603 proposed STL-based approach can provide more promising and reliable results for
604 Facial Expression Recognition. By eliminating the need for many labels in the early
605 stages of training, the STL approach has proven to be efficient and cost-effective while
606 achieving competitive results in terms of accuracy and generalization.

607 5 Conclusion

608 We have proposed a new method based on STL applied to the FER problem. The
609 primary advancements within our self-taught learning approach, as compared to the
610 existing literature, can be summarized as (a) the REL algorithm, which can generate
611 a pool with diverse representations using an unsupervised dataset, and (b) a robust
612 performance evaluation, drawing a direct comparison between our suggested unsuper-
613 vised FER solution and supervised methods within identical conditions, this entails
614 utilizing the same experimental protocol and addressing the challenges posed by a
615 genuine and challenging task such as FER.

616 As future work, we plan to investigate different directions. First, we intend to
617 evaluate an alternative scheme for the loss function employed while learning the unsu-
618 pervised set of latent representations, aiming at amplifying the diversity. Second, we

619 plan to extend the evaluation of the proposed REL algorithm to cover different target
620 problems, exploring its applicability and effectiveness in other contexts. In addition,
621 we plan to consider the investigation of transformers and large vision models, espe-
622 cially to evaluate whether these architectures provide additional performance gains
623 compared to those explored in our work. These initiatives aim to improve further the
624 scope and robustness of our contributions.

625 Acknowledgement

626 We would like to thank the *Conselho Nacional de Desenvolvimento Científico e Tec-*
627 *nológico* (CNPq, grants 306878/2022-4, 406030/2023-5 and 441610/2023-4) and *Coor-*
628 *denação de Aperfeiçoamento de Pessoal de Nível Superior* (CAPES) for co-financing
629 this research.

630 6 Declarations

631 This section summarizes some important statements as follows:

- 632 • **Funding and/or Conflicts of Interests/Competing Interests:** The authors
633 certify that they have no conflict of interest.
- 634 • **Data Availability:** The datasets used in the experiments are available to the
635 scientific community upon request and sign of a proper responsibility agreement.
- 636 • **Authorship Confirmation:** all authors have participated in the conception,
637 design, analysis, and interpretation of the data, drafting the article or revising it,
638 and approving the final version.

639 References

- 640 [1] Yeung, M.K.: A systematic review and meta-analysis of facial emotion recogni-
641 tion in autism spectrum disorder: The specificity of deficits and the role of task
642 characteristics. *Neuroscience & Biobehavioral Reviews* **133**, 104518 (2022)
- 643 [2] Vehlen, A., Kellner, A., Normann, C., Heinrichs, M., Domes, G.: Reduced eye
644 gaze during facial emotion recognition in chronic depression: Effects of intranasal
645 oxytocin. *Journal of Psychiatric Research* **159**, 50–56 (2023)
- 646 [3] Qiao, Y., Zeng, K., Xu, L., Yin, X.: A smartphone-based driver fatigue detection
647 using fusion of multiple real-time facial features. In: 2016 13th IEEE Annual Con-
648 sumer Communications & Networking Conference (CCNC), pp. 230–235 (2016).
649 IEEE
- 650 [4] Li, S., Deng, W.: Deep facial expression recognition: A survey. *IEEE transactions*
651 *on affective computing* **13**(3), 1195–1215 (2020)
- 652 [5] Rathour, N., Singh, R., Gehlot, A., Akram, S.V., Thakur, A.K., Kumar, A.:
653 The decadal perspective of facial emotion processing and recognition: A survey.
654 *Displays*, 102330 (2022)

- 655 [6] Ganaie, M.A., Hu, M., Malik, A., Tanveer, M., Suganthan, P.: Ensemble deep
656 learning: A review. *Engineering Applications of Artificial Intelligence* **115**, 105151
657 (2022)
- 658 [7] Raina, R., Battle, A., Lee, H., Packer, B., Ng, A.Y.: Self-taught learning: trans-
659 fer learning from unlabeled data. In: *Proceedings of the 24th International*
660 *Conference on Machine Learning*, pp. 759–766 (2007)
- 661 [8] Bengio, Y., Bastien, F., Bergeron, A., Boulanger-Lewandowski, N., Chherawala,
662 Y., Cisse, M., Côté, M., Erhan, D., Eustache, J., Glorot, X., *et al.*: Deep self-
663 taught learning for handwritten character recognition. In: *NIPS* 2010 Deep*
664 *Learning and Unsupervised Feature Learning Workshop* (2010)
- 665 [9] Ojala, T., Pietikäinen, M., Harwood, D.: A comparative study of texture measures
666 with classification based on featured distributions. *Pattern recognition* **29**(1),
667 51–59 (1996)
- 668 [10] Zavaschi, T.H., Britto Jr, A.S., Oliveira, L.E., Koerich, A.L.: Fusion of feature sets
669 and classifiers for facial expression recognition. *Expert Systems with Applications*
670 **40**(2), 646–655 (2013)
- 671 [11] Lowe, D.G.: Distinctive image features from scale-invariant keypoints. *Internation-*
672 *al journal of computer vision* **60**, 91–110 (2004)
- 673 [12] Haralick, R.M., Shanmugam, K., Dinstein, I.H.: Textural features for image clas-
674 sification. *IEEE Transactions on systems, man, and cybernetics* (6), 610–621
675 (1973)
- 676 [13] Bay, H., Ess, A., Tuytelaars, T., Van Gool, L.: Speeded-up robust features (surf).
677 *Computer vision and image understanding* **110**(3), 346–359 (2008)
- 678 [14] Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection.
679 In: *2005 IEEE Computer Society Conference on Computer Vision and Pattern*
680 *Recognition (CVPR'05)*, vol. 1, pp. 886–893 (2005). Ieee
- 681 [15] Canal, F.Z., Müller, T.R., Matias, J.C., Scotton, G.G., Sa Junior, A.R., Pozzebon,
682 E., Sobieranski, A.C.: A survey on facial emotion recognition techniques: A state-
683 of-the-art literature review. *Information Sciences* **582**, 593–617 (2022)
- 684 [16] Feng, S., Yu, H., Duarte, M.F.: Autoencoder based sample selection for self-taught
685 learning. *Knowledge-Based Systems* **192**, 105343 (2020)
- 686 [17] Allognon, S.O.C., Britto, A.d.S., Koerich, A.L.: Continuous emotion recogni-
687 tion via deep convolutional autoencoder and support vector regressor. In: *2020*
688 *International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8 (2020).
689 IEEE

- 690 [18] Delazeri, B.R., Vera, L.L., Barddal, J.P., Koerich, A.L., *et al.*: Evaluation of
691 self-taught learning-based representations for facial emotion recognition. In: 2022
692 International Joint Conference on Neural Networks (IJCNN), pp. 1–8 (2022).
693 IEEE
- 694 [19] Renda, A., Barsacchi, M., Bechini, A., Marcelloni, F.: Comparing ensemble strate-
695 gies for deep learning: An application to facial expression recognition. *Expert*
696 *Systems with Applications* **136**, 1–11 (2019)
- 697 [20] Wen, G., Hou, Z., Li, H., Li, D., Jiang, L., Xun, E.: Ensemble of deep neural
698 networks with probability-based fusion for facial expression recognition. *Cognitive*
699 *Computation* **9**(5), 597–610 (2017)
- 700 [21] Li, R., Ren, C., Zhang, X., Hu, B.: A novel ensemble learning method using
701 multiple objective particle swarm optimization for subject-independent eeg-based
702 emotion recognition. *Computers in biology and medicine* **140**, 105080 (2022)
- 703 [22] Dhankhar, P.: Resnet-50 and vgg-16 for recognizing facial emotions. *International*
704 *Journal of Innovations in Engineering and Technology (IJJET)* **13**(4), 126–130
705 (2019)
- 706 [23] Chowdary, M.K., Nguyen, T.N., Hemanth, D.J.: Deep learning-based facial emo-
707 tion recognition for human–computer interaction applications. *Neural Computing*
708 *and Applications* **35**(32), 23311–23328 (2023)
- 709 [24] Akhand, M., Roy, S., Siddique, N., Kamal, M.A.S., Shimamura, T.: Facial emotion
710 recognition using transfer learning in the deep cnn. *Electronics* **10**(9), 1036 (2021)
- 711 [25] Lee, H., Grosse, R., Ranganath, R., Ng, A.Y.: Unsupervised learning of hierar-
712 chical representations with convolutional deep belief networks. *Communications*
713 *of the ACM* **54**(10), 95–103 (2011)
- 714 [26] Markov, K., Matsui, T.: Music genre classification using self-taught learning via
715 sparse coding. In: 2012 IEEE International Conference on Acoustics, Speech and
716 Signal Processing (ICASSP), pp. 1929–1932 (2012). IEEE
- 717 [27] Hu, H., Phan, N., Chun, S.A., Geller, J., Vo, H., Ye, X., Jin, R., Ding, K., Kenne,
718 D., Dou, D.: An insight analysis and detection of drug-abuse risk behavior on
719 twitter with self-taught deep learning. *Computational Social Networks* **6**(1), 1–19
720 (2019)
- 721 [28] Qureshi, A.S., Khan, A., Shamim, N., Durad, M.H.: Intrusion detection using
722 deep sparse auto-encoder and self-taught learning. *Neural Computing and*
723 *Applications* **32**, 3135–3147 (2020)
- 724 [29] Li, S., Li, K., Fu, Y.: Self-taught low-rank coding for visual learning. *IEEE*
725 *transactions on neural networks and learning systems* **29**(3), 645–656 (2017)

- 726 [30] Ramamurthy, S.R., Ghosh, I., Gangopadhyay, A., Galik, E., Roy, N.: Star-lite:
727 A light-weight scalable self-taught learning framework for older adults' activity
728 recognition. *Pervasive and Mobile Computing* **87**, 101698 (2022)
- 729 [31] Germani, E., Fromont, E., Maumet, C.: On the benefits of self-taught learning
730 for brain decoding. *GigaScience* **12**, 029 (2023)
- 731 [32] He, P., Jia, P., Qiao, S., Duan, S.: Self-taught learning based on sparse
732 autoencoder for e-nose in wound infection detection. *Sensors* **17**(10), 2279 (2017)
- 733 [33] Liu, L., Wei, W., Chow, K.-H., Loper, M., Gursoy, E., Truex, S., Wu, Y.: Deep
734 neural network ensembles against deception: Ensemble diversity, accuracy and
735 robustness. In: 2019 IEEE 16th International Conference on Mobile Ad Hoc and
736 Sensor Systems (MASS), pp. 274–282 (2019). IEEE
- 737 [34] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,
738 Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., *et al.*: Scikit-learn: Machine
739 learning in python. *Journal of machine learning research* **12**(Oct), 2825–2830
740 (2011)
- 741 [35] Lyons, M., Akamatsu, S., Kamachi, M., Gyoba, J.: Coding facial expressions
742 with gabor wavelets. In: Proceedings Third IEEE International Conference on
743 Automatic Face and Gesture Recognition, pp. 200–205 (1998). IEEE
- 744 [36] Lucey, P., Cohn, J.F., Kanade, T., Saragih, J., Ambadar, Z., Matthews, I.:
745 The extended cohn-kanade dataset (ck+): A complete dataset for action unit
746 and emotion-specified expression. In: 2010 Ieee Computer Society Conference on
747 Computer Vision and Pattern Recognition-workshops, pp. 94–101 (2010). IEEE
- 748 [37] Viola, P., Jones, M.J.: Robust real-time face detection. *International journal of*
749 *computer vision* **57**, 137–154 (2004)
- 750 [38] Doi, E., Inui, T., Lee, T.-W., Wachtler, T., Sejnowski, T.J.: Spatiochromatic
751 receptive field properties derived from information-theoretic analyses of cone
752 mosaic responses to natural scenes. *Neural computation* **15**(2), 397–417 (2003)
- 753 [39] Uetz, R., Behnke, S.: Large-scale object recognition with cuda-accelerated hier-
754 archical neural networks. In: 2009 IEEE International Conference on Intelligent
755 Computing and Intelligent Systems, vol. 1, pp. 536–541 (2009). IEEE
- 756 [40] Huang, G.B., Mattar, M., Berg, T., Learned-Miller, E.: Labeled faces in the
757 wild: A database for studying face recognition in unconstrained environments. In:
758 Workshop on Faces in 'Real-Life' Images: Detection, Alignment, and Recognition
759 (2008)
- 760 [41] Alzubaidi, L., Zhang, J., Humaidi, A.J., Al-Dujaili, A., Duan, Y., Al-Shamma,
761 O., Santamaría, J., Fadhel, M.A., Al-Amidie, M., Farhan, L.: Review of deep

- 762 learning: Concepts, cnn architectures, challenges, applications, future directions.
763 *Journal of big Data* **8**, 1–74 (2021)
- 764 [42] Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep
765 convolutional neural networks. *Advances in neural information processing systems*
766 **25** (2012)
- 767 [43] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale
768 image recognition. *arXiv preprint arXiv:1409.1556* (2014)
- 769 [44] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D.,
770 Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: *Proceedings*
771 *of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–9
772 (2015)
- 773 [45] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recogni-
774 tion. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern*
775 *Recognition*, pp. 770–778 (2016)
- 776 [46] Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected
777 convolutional networks. In: *Proceedings of the IEEE Conference on Computer*
778 *Vision and Pattern Recognition*, pp. 4700–4708 (2017)
- 779 [47] Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neu-
780 ral networks. In: *International Conference on Machine Learning*, pp. 6105–6114
781 (2019). PMLR
- 782 [48] Xu, M., Cheng, W., Zhao, Q., Ma, L., Xu, F.: Facial expression recognition based
783 on transfer learning from deep convolutional networks. In: *2015 11th International*
784 *Conference on Natural Computation (ICNC)*, pp. 702–708 (2015). IEEE
- 785 [49] Zhang, J., Li, W., Ogunbona, P., Xu, D.: Recent advances in transfer learn-
786 ing for cross-dataset visual recognition: A problem-oriented perspective. *ACM*
787 *Computing Surveys (CSUR)* **52**(1), 1–38 (2019)
- 788 [50] Peng, M., Wu, Z., Zhang, Z., Chen, T.: From macro to micro expression recogni-
789 tion: Deep learning on small datasets using transfer learning. In: *2018 13th IEEE*
790 *International Conference on Automatic Face & Gesture Recognition (FG 2018)*,
791 pp. 657–661 (2018). IEEE
- 792 [51] Shao, J., Qian, Y.: Three convolutional neural network models for facial expres-
793 sion recognition in the wild. *Neurocomputing* **355**, 82–92 (2019)
- 794 [52] Kola, D.G.R., Samayamantula, S.K.: Facial expression recognition using singular
795 values and wavelet-based lgc-hd operator. *IET Biometrics* **10**(2), 207–218 (2021)

- 796 [53] Kartheek, M.N., Prasad, M.V., Bhukya, R.: Radial mesh pattern: a hand-
797 crafted feature descriptor for facial expression recognition. *Journal of Ambient*
798 *Intelligence and Humanized Computing* **14**(3), 1619–1631 (2023)
- 799 [54] Mandal, M., Verma, M., Mathur, S., Vipparthi, S.K., Murala, S., Kranthi Kumar,
800 D.: Regional adaptive affinitive patterns (radap) with logical operators for facial
801 expression recognition. *IET Image Processing* **13**(5), 850–861 (2019)
- 802 [55] Du, L., Hu, H.: Weighted patch-based manifold regularization dictionary pair
803 learning model for facial expression recognition using iterative optimization
804 classification strategy. *Computer Vision and Image Understanding* **186**, 13–24
805 (2019)
- 806 [56] Wu, B.-F., Lin, C.-H.: Adaptive feature mapping for customizing deep learning
807 based facial expression recognition model. *IEEE access* **6**, 12451–12461 (2018)
- 808 [57] Lee, S.H., Baddar, W.J., Ro, Y.M.: Collaborative expression representation using
809 peak expression and intra class variation face images for practical subject-
810 independent emotion recognition in videos. *Pattern Recognition* **54**, 52–67
811 (2016)